

Asymptotic properties of the number of matching coalescent histories for caterpillar-like families of species trees

Filippo Disanto*
Noah A. Rosenberg

Department of Biology, Stanford University, Stanford, CA 94305 USA

March 20, 2015

Abstract

Coalescent histories provide lists of species tree branches on which gene tree coalescences can take place, and their enumerative properties assist in understanding the computational complexity of calculations central in the study of gene trees and species trees. Here, we solve an enumerative problem left open by Rosenberg (*IEEE/ACM Transactions on Computational Biology and Bioinformatics* 10: 1253-1262, 2013) concerning the number of coalescent histories for gene trees and species trees with a matching labeled topology that belongs to a generic caterpillar-like family. By bringing a generating function approach to the study of coalescent histories, we prove that for any caterpillar-like family with seed tree t , the sequence $(h_n)_{n \geq 0}$ describing the number of matching coalescent histories of the n th tree of the family grows asymptotically as a constant multiple of the Catalan numbers. Thus, $h_n \sim \beta_t c_n$, where the asymptotic constant $\beta_t > 0$ depends on the shape of the seed tree t . The result extends a claim demonstrated only for seed trees with at most 8 taxa to arbitrary seed trees, expanding the set of cases for which detailed enumerative properties of coalescent histories can be determined. We introduce a procedure that computes from t the constant β_t as well as the algebraic expression for the generating function of the sequence $(h_n)_{n \geq 0}$.

1 Introduction

Coalescent histories, mathematical structures representing combinatorially distinct ways in which a given gene tree can coalesce along the branches of a given species tree, are important in a variety of phylogenetic problems [5, 13, 14]. They arise most prominently in characterizing the set of objects over which a sum is performed in a fundamental calculation for inference of species trees from information on multiple genetic loci, the evaluation of gene tree probabilities conditional on species trees [4].

Because of the appearance of coalescent histories in sets over which sums are computed, as well as in state spaces of certain phylogenetic Markov chains [6, 9, 10], solutions to enumerative problems involving coalescent histories contribute to an understanding of the computational complexity of phylogenetic calculations. A recursion for the number of coalescent histories for a given gene tree and species tree has been established [12], and several studies have reported exact numerical results and closed-form expressions for the number of coalescent histories for small trees and for specific types of trees of arbitrarily large size [3, 4, 5, 12, 13, 14, 16]. The latter computations have proceeded both by solving or deploying the recursion in specific cases [12, 13, 14, 16], as well as by identifying correspondences between coalescent histories and other combinatorial structures for which enumerative results have already been established [3, 4, 5].

*Corresponding author. Email: fdisanto@stanford.edu.

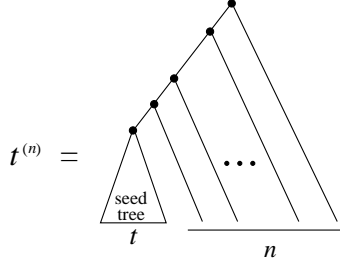


Figure 1: A caterpillar-like family of species trees $(t^{(n)})_{n \geq 0}$. For a seed tree t , by adding $n \geq 0$ branches each with 1 leaf, we obtain the n th tree of the family, $t^{(n)}$. If t has 2 taxa, then $(t^{(n)})_{n \geq 0}$ is simply the caterpillar family.

One class of gene trees and species trees of particular interest for enumeration of coalescent histories is the *caterpillar-like families*, trees that have a caterpillar shape, except that the caterpillar subtree with r taxa is replaced by a subtree of size r that is not necessarily a caterpillar subtree (Fig. 1). For the simplest caterpillar-like family, the set of caterpillar trees themselves, if the gene tree and species tree have the same caterpillar labeled topology with n taxa, then the number of coalescent histories is a Catalan number,

$$c_{n-1} = \frac{1}{n} \binom{2n-2}{n-1}. \quad (1)$$

For T_r -caterpillar-like families, in which the r -taxon subtree of an n -taxon caterpillar species tree is replaced by an r -taxon subtree T_r (Fig. 1), by employing the recursion method, Rosenberg [13] obtained the exact number of coalescent histories for all n , for each T_r with $r \leq 8$, in the case that the gene tree and species tree have the same labeled topology. Rosenberg [13] argued that in each of these cases, as $n \rightarrow \infty$, the number of coalescent histories is asymptotic to a constant multiple of the Catalan numbers. A proof of this result has been presented in full for each case with $r \leq 5$ [3, 12, 13], and by computer algebra for cases with $r = 6, 7$, and 8 [13].

Here, using a substantially different approach that brings to studies of coalescent histories the methods of analytic combinatorics, we produce an enumeration result that covers caterpillar-like families in general. We show that the result of [13] applies to all caterpillar-like families, not only those for which T_r has $r \leq 8$. That is, we demonstrate that for any T_r , as $n \rightarrow \infty$, the number of coalescent histories in the T_r -caterpillar-like family is asymptotic to a constant multiple of the Catalan numbers. We describe a method for computing the constant and provide a symbolic tool for performing the computation. Finally, we discuss the results in terms of their impact in mathematical phylogenetics.

2 Preliminaries

2.1 Species trees and coalescent histories

We consider binary rooted leaf-labeled species trees, taking a single arbitrary labeling (without loss of generality) to represent a given unlabeled species tree topology. We consider an arbitrarily labeled species tree and its unlabeled tree interchangeably, treating the labeling as implicit.

We examine coalescent histories for the case in which gene trees and species trees have the same labeled topology t , terming a coalescent history in this case a *matching coalescent history*. To be a matching coalescent history, a mapping h from the internal nodes of t (viewed as the gene tree) to the branches of t (viewed as the species tree) must satisfy two conditions: (a) for each leaf x in t , if x descends from node k in t , then x descends from branch $h(k)$ in t ; (b) for each pair of internal nodes k_1 and k_2 in t , if k_2 descends from k_1 in t , then branch $h(k_2)$ descends from or coincides with branch $h(k_1)$ in t . The definition of matching coalescent histories is illustrated in Figure 2. We henceforth consider only matching coalescent histories, treating “matching” as implicit in references to coalescent histories; we also refer simply to *histories* for short.

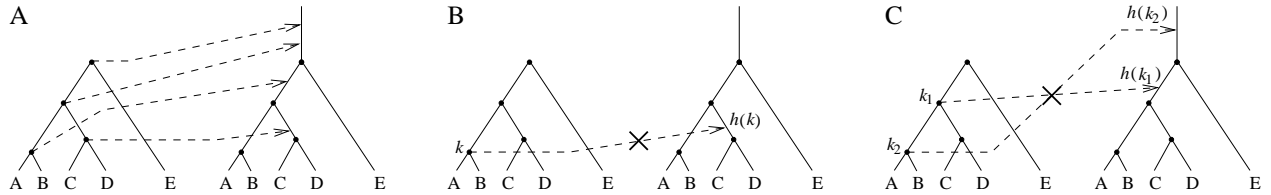


Figure 2: Matching coalescent histories. (A) A matching coalescent history. (B) A mapping from the internal nodes of a tree to its branches that does not satisfy condition (a). Leaf B is descended from node k but does not descend from branch $h(k)$. (C) A mapping from the internal nodes of a tree to its internal branches that does not satisfy condition (b). Node k_2 is descended from node k_1 , but branch $h(k_2)$ is strictly ancestral to branch $h(k_1)$.

2.2 Caterpillar-like families of species trees

For a binary species tree t with at least 2 taxa, we denote by $(t^{(n)})_{n \geq 0}$ the caterpillar-like family generated by the seed tree t . This family is recursively defined by taking $t^{(0)} = t$ and letting $t^{(n+1)}$ be the tree obtained by appending $t^{(n)}$ and a single leaf to a same root (Fig. 1).

Our interest is in the number of matching coalescent histories of $t^{(n)}$ for $n \geq 0$, a quantity we denote by $h_n(t)$ or simply h_n . We note that whereas [13] indexed trees by their numbers of taxa, here n represents the number of taxa appended above the root of the seed tree, so that if seed tree t has $|t|$ taxa, then $|t| + n$ gives the number of taxa in $t^{(n)}$.

2.3 Principles of analytic combinatorics

We rely on techniques of analytic combinatorics [7] to obtain our enumerative results, and recall several key points. In general, an integer sequence $(a_n)_{n \geq 0}$ can be associated with a formal power series $A(z) = \sum_{n=0}^{\infty} a_n z^n$, also termed the *generating function* of the integers a_n . Considering z as a complex variable, typically in a neighborhood of 0, features of the function $A(z)$ are related to the growth of the coefficients a_n .

More precisely, generating functions, considered as complex functions, enable analyses of the asymptotic growth of the associated integer sequences through the analysis of their singularities in the complex plane. In particular, under suitable conditions, there exists a general correspondence between the singular expansion of a generating function $A(z)$ near its dominant singularities—those nearest the origin—and the asymptotic behavior of the associated coefficients a_n (Chapter VI of [7]). We make use of generating functions that near their unique dominant singularity can be described by means of the square root function, and for which theorems on singularity analysis of generating functions [7] consequently apply.

2.4 Catalan numbers

The Catalan sequence appears often in combinatorics [7, 8, 15] and features prominently in our analysis. Rewriting eq. (1) with index n rather than $n - 1$,

$$c_n = \frac{1}{n+1} \binom{2n}{n}. \quad (2)$$

The associated generating function is

$$C(z) = \sum_{n=0}^{\infty} c_n z^n = \frac{1 - \sqrt{1 - 4z}}{2z}. \quad (3)$$

By definition, if $[z^n]f(z)$ denotes the n th term in the power series expansion of $f(z)$ at $z = 0$, we have

$$c_n = [z^n]C(z) = \frac{1}{2}[z^{n+1}](1 - \sqrt{1 - 4z}) = \frac{1}{2}[z^{n+1}](-\sqrt{1 - 4z}). \quad (4)$$

Asymptotically, applying Stirling's approximation $n! \sim \sqrt{2\pi n}(n/e)^n$ to eq. (2), the Catalan sequence satisfies

$$c_n \sim \frac{4^n}{n^{3/2}\sqrt{\pi}}. \quad (5)$$

3 The number of matching coalescent histories for caterpillar-like families

Our goal is to produce a procedure that evaluates the number of coalescent histories $h_n(t)$ for matching gene trees and species trees in the caterpillar-like family that begins with seed tree t , and moreover, to show that

$$h_n(t) \sim \beta_t c_n, \quad (6)$$

where the multiplier $\beta_t > 0$ for the Catalan sequence is a constant depending on t . In other words, we wish to demonstrate that as $n \rightarrow \infty$, the ratio h_n/c_n converges to a constant $\beta_t > 0$ that depends on the seed tree t .

First, in Section 3.1, we determine a lower bound for the number of matching coalescent histories of the n th tree $t^{(n)}$ of the caterpillar-like family with seed tree t . Next, in Section 3.2, we introduce a concept of *m-rooted histories* of a species tree $t^{(n)}$. The section provides an iterative construction of the rooted histories of $t^{(n+1)}$ from those of $t^{(n)}$, describing the construction by means of a convenient labeling scheme. We follow a commonly used combinatorial enumeration strategy [1, 2] that determines a recursive succession rule for successive collections of objects in a sequence and then uses this rule to compute a generating function. In Section 3.3, we use the iterative construction to produce a bivariate generating function whose coefficients $h_{n,m}$ are the numbers of *m*-rooted histories for trees $t^{(n)}$. We next obtain the generating function for the integer sequence $(h_n)_{n \geq 0}$ describing the number of matching coalescent histories for the $t^{(n)}$. Finally, using the lower bound from Section 3.1, in Section 3.4, we apply methods of analytic combinatorics to study the asymptotic behavior of h_n .

3.1 Lower bound for h_n

To produce a lower bound for h_n , we first define V as the tree with 2 taxa. Recalling that we index trees so that the number of taxa in a tree is n more than the number of taxa in the seed tree, we have [3, 12, 13]

$$h_n(V) = c_{n+1}.$$

A constructive procedure, illustrated in Figure 3, shows that for any seed tree t with $|t| \geq 2$,

$$h_n(t) \geq h_n(V) = c_{n+1}. \quad (7)$$

For a seed tree t , we can superimpose V on t so that the root r_V of V matches the root r_t of t (Fig. 3B). The two leaves of V are identified with two of the leaves of t , one on each side of the root of t . Generating caterpillar-like families by adding n single branches separately to V and to t , the superposition of V on t extends, so that $V^{(n)}$ is superimposed on $t^{(n)}$ (Fig. 3C). The n caterpillar branches of $t^{(n)}$ and $V^{(n)}$ then correspond.

Each matching coalescent history h of $t^{(n)}$ determines a corresponding matching coalescent history h' of $V^{(n)}$ by considering the restriction of the history h to the set of internal nodes of $t^{(n)}$ that correspond to internal nodes of $V^{(n)}$ (Fig. 3D). Thus, for any given seed tree t , the number of matching coalescent histories of $t^{(n)}$ is greater than or equal to the number of matching coalescent histories of $V^{(n)}$. In symbols, we have eq. (7).

3.2 Iterative generation of rooted histories

This section describes the iterative procedure that for a seed tree t eventually enables us to determine a formula for h_n . First, in Section 3.2.1, we discuss *m*-rooted histories, which extend the concept of matching coalescent histories, introducing an additional parameter m . Next, in Section 3.2.2, we examine the relationship between

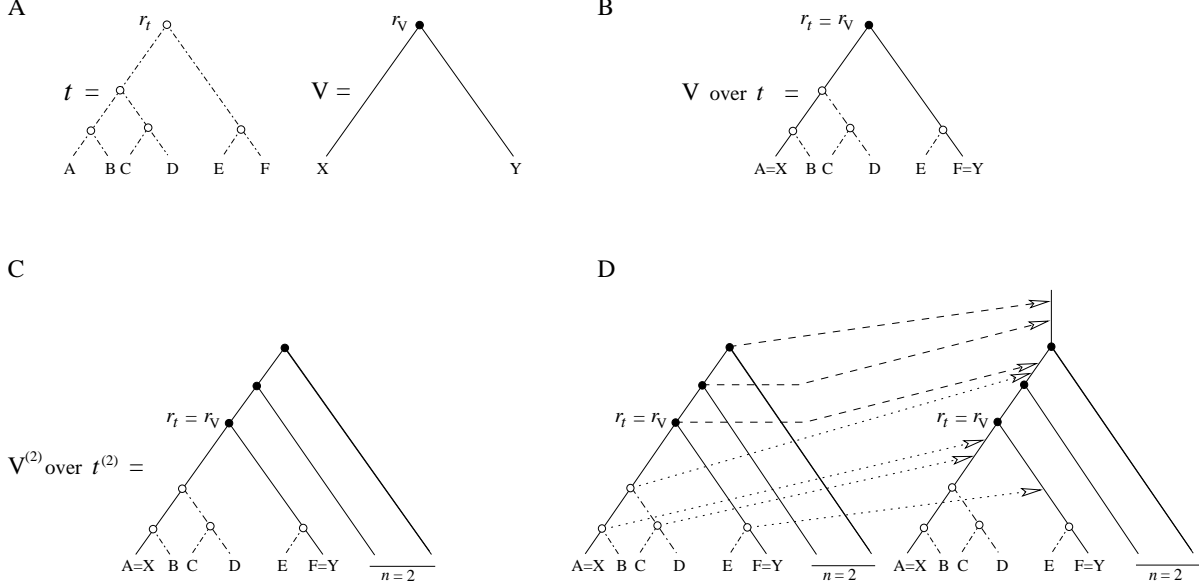


Figure 3: Superposition of the caterpillar tree family on a caterpillar-like tree family with arbitrary seed tree of size $|t| \geq 2$. (A) A seed tree t and the seed tree V for the caterpillar family. (B) Superposition of V on t , so that the roots r_V and r_t overlap. (C) Superposition of $V^{(2)}$ (shaded internal nodes) on $t^{(2)}$ (shaded and unshaded nodes). The $n = 2$ caterpillar branches in $V^{(2)}$ and $t^{(2)}$ overlap, and r_V still matches r_t . (D) A matching coalescent history of $t^{(2)}$ (dashed and dotted arrows) determines a matching coalescent history of $V^{(2)}$ (dashed arrows) by ignoring arrows from the unshaded nodes.

rooted histories and the *extended coalescent histories* of [12], importing results on extended coalescent histories into the more convenient framework of rooted histories. We expand our goal of enumerating matching coalescent histories for $t^{(n)}$, considering a more general problem of enumerating for $m \geq 1$ the m -rooted histories of $t^{(n)}$.

In Section 3.2.3, we define an operator Ω for constructing the rooted histories of $t^{(n+1)}$ from the rooted histories of $t^{(n)}$. Next, in Section 3.2.4, we introduce a labeling scheme that in Section 3.2.5 enables us to switch from counting rooted histories to counting multisets of labels. At the end of Section 3.2, we will have converted our enumeration problem into an enumeration that is more convenient for constructing a generating function.

3.2.1 m -rooted histories

Consider a tree t with $|t| \geq 2$, and suppose that the branch above the root of t (the *root-branch*) is divided into infinitely many components. A matching coalescent history mapping the internal nodes of t onto the branches of t is said to be m -rooted for $m \geq 1$ if the root of t is mapped *exactly* onto the m th component of the root (Fig. 4). It is said to be *rooted* if it is m -rooted for some m . Branches are numbered so that branch $m = 1$ is immediately above the root node, and m is greater for components that are farther from the root.

For a rooted history h of a tree t , $m = m(h)$ denotes the component of the root-branch of t that receives the image of the root of t . $H_{n,m}(t)$ denotes the set of m -rooted histories of $t^{(n)}$, and $H_n(t) = \bigcup_{m=1}^{\infty} H_{n,m}(t)$ denotes the set of its rooted histories. The number of m -rooted histories of $t^{(n)}$ is $h_{n,m} = |H_{n,m}|$, and the number of 1-rooted histories $h_n = h_{n,1}$ is also the number of matching coalescent histories. Enumeration of the matching coalescent histories of $t^{(n)}$ is equivalent to enumeration of the 1-rooted histories of $t^{(n)}$.

3.2.2 Rooted histories and extended histories

Rooted histories are closely related to *extended coalescent histories*, as defined by [12]. We use this relationship to study properties of rooted histories. Rosenberg [12] defined the set of k -*extended coalescent histories* of a tree

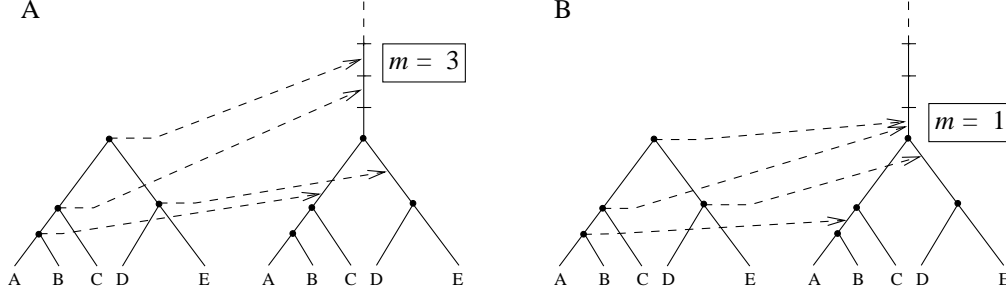


Figure 4: Rooted histories of a tree. (A) A 3-rooted history. The root-branch is divided into infinitely many components. The third component receives the image of the root. (B) A 1-rooted history. The number of 1-rooted histories corresponds to the number of matching coalescent histories of the tree.

t with $|t| \geq 1$ for integers $k \geq 1$; we also consider $k = 0$ by setting the number of 0-extended histories to 0.

A k -extended history is defined as a coalescent history for a species tree whose root-branch is divided into exactly $k \geq 0$ parts. In other words, the root-branch has exactly $k \geq 0$ possible components onto which a k -extended history can map the gene tree root. Here we consider matching k -extended histories, so that the internal nodes of a tree t are mapped to the branches of t and its k components above the root. For convenience, we refer to extended histories by the index k , reserving the index m for rooted histories.

By the definitions of k -extended and m -rooted histories, for each $k \geq 0$, the set of k -extended histories of a tree is exactly the set of all m -rooted histories with $1 \leq m \leq k$. Therefore, for a tree t with at least 2 leaves, if we label by $e_{t,k}$ its number of k -extended histories, then for each $m \geq 1$ the number of m -rooted histories of t is

$$h_{0,m} = e_{t,m} - e_{t,m-1}. \quad (8)$$

Note that for $m = 1$, we explicitly use in eq. (8) the fact that $e_{t,0}$ is defined and equal to 0. In addition to setting $e_{t,0} = 0$ for any tree t , as in [12] we also set $e_{t,k} = 1$ for all $k \geq 1$ in the case that t has exactly 1 leaf.

Suppose $|t| \geq 1$ and $k \geq 0$. Denote by t_L and t_R the left and right subtrees of the root of t . We can compute $e_{t,k}$ recursively as in Theorem 3.1 of [12]:

$$e_{t,k} = \begin{cases} 0 & \text{if } |t| \geq 1 \text{ and } k = 0 \\ 1 & \text{if } |t| = 1 \text{ and } k \geq 1 \\ \sum_{i=1}^k e_{t_L,i+1} e_{t_R,i+1} & \text{if } |t| \geq 2 \text{ and } k \geq 1. \end{cases} \quad (9)$$

As was already observed in the remarks following Corollary 3.2 of [12], by eq. (9), for any tree t with $|t| \geq 1$, for positive integers $k \geq 1$, the function $f(k) = e_{t,k}$ is a polynomial in k . With our extension to permit $k = 0$, we can extend this fact to $k \geq 0$ for $|t| \geq 2$: for any tree t with $|t| \geq 2$, and for $k \geq 0$, we claim that the function $f(k) = e_{t,k}$ is a polynomial in k . Note that in allowing $k = 0$, we claim $e_{t,k}$ is a polynomial in k only for $|t| \geq 2$; for $|t| = 1$, $e_{t,k}$ is not a polynomial in k because $e_{t,0} = 0$ and $e_{t,k} = 1$ for $k \geq 1$.

To prove the claim, fix t with $|t| \geq 2$ and consider the variable k over domain $[1, \infty)$. We demonstrate that $f(k)$ is a polynomial in k for domain $[0, \infty)$ by showing that the closed-form polynomial for $f(k)$ has a factor of k , so that our choice $e_{t,0} = 0$ in eq. (9) is compatible with the polynomial expression valid for $k \geq 1$.

Observe that for $i \geq 1$, $e_{t_L,i}$ and $e_{t_R,i}$ are polynomials in i , say $P_{t_L}(i)$ and $P_{t_R}(i)$. Replacing the terms $e_{t_L,i+1}$ and $e_{t_R,i+1}$ that appear in the recursion in eq. (9) by the two polynomials $P_{t_L}(i+1)$ and $P_{t_R}(i+1)$, we obtain

$$\sum_{i=1}^k e_{t_L,i+1} e_{t_R,i+1} = \sum_{i=1}^k P_{t_L}(i+1) P_{t_R}(i+1) = \sum_{i=1}^k P'(i), \quad (10)$$

where $P'(i)$ denotes a polynomial in i that results from the product of $P_{t_L}(i+1)$ and $P_{t_R}(i+1)$. By Faulhaber's formula for sums of powers of integers, symbolic sums of the form $\sum_{i=1}^k i^p$ for a fixed integer $p \geq 0$ are polynomials

containing a factor of k in their closed forms (Section 6.5 of [8])—for example, $\sum_{i=1}^k i^3 = k^2(k+1)^2/4$. Thus, because the polynomial $P'(i)$ is a linear combination of terms of the form i^p , the closed-form expression for the sum $\sum_{i=1}^k P'(i)$ appearing in eq. (10) also has a factor of k . It therefore has a value of 0 at $k = 0$.

Functions $e_{t,k}$ for trees t with $1 \leq |t| \leq 9$ and $k \geq 1$ appear in Tables 1-4 of [12]. For $|t| \geq 2$, as we have shown, these example polynomials are divisible by the variable representing the number of components of the root-branch. By eq. (8), we immediately obtain the following result.

Proposition 1 *For any tree t with $|t| \geq 2$ and for $m \geq 1$, the number $h_{0,m}$ of m -rooted histories of t is a polynomial in m that can be computed by the difference in eq. (8) using $e_{t,k}$ as in eq. (9).*

As an example of Proposition 1, consider the tree $t = ((A, B), (C, D))$, identifying this arbitrary labeling with the unlabeled tree $((())())$. By applying the recursive procedure in eq. (9), we find that for $k \geq 0$, the number of k -extended coalescent histories for t is $e_{t,k} = \frac{1}{6}k(2k^2 + 9k + 13)$ [12]. The difference eq. (8) yields that for $m \geq 1$ the number of m -rooted histories of t is $h_{0,m} = e_{t,m} - e_{t,m-1} = m^2 + 2m + 1$.

3.2.3 Generating rooted histories of $t^{(n+1)}$ from rooted histories of $t^{(n)}$

This section introduces an operator Ω that generates the rooted histories of $t^{(n+1)}$ from those of $t^{(n)}$. For each rooted history h' of $t^{(n+1)}$, there exists exactly one rooted history h of $t^{(n)}$ with $h' \in \Omega(h)$. Recalling the definitions of the sets $H_{n,m}(t)$ and $H_n(t)$ of m -rooted and rooted histories of $t^{(n)}$, we define Ω as follows.

Definition. Let $\mathcal{P}(X) = \{x : x \subseteq X\}$ denote the power set of set X , and fix tree t . The operator Ω is a function

$$\Omega : H_n(t) \rightarrow \mathcal{P}(H_{n+1}(t)),$$

where for a given rooted history $h \in H_n(t)$, $\Omega(h)$ is the set of rooted histories $h' \in H_{n+1}(t)$ for which the restriction of h' to $t^{(n+1)}$ excluding its most basal caterpillar branch coincides with the rooted history h of $t^{(n)}$.

Denote by b_1, b_2, \dots, b_{n+1} the caterpillar branches in $t^{(n+1)}$, from the least basal b_1 to the most basal b_{n+1} (Fig. 5). Upon removal of the most basal caterpillar branch b_{n+1} from $t^{(n+1)}$, the root of $t^{(n+1)}$ —to which branch b_{n+1} is attached—is replaced by a demarcation between the first and second components of the root-branch of $t^{(n)}$. For instance, in Fig. 5A, starting from tree $t = ((A, B), (C, D))$, we consider h''' , a 3-rooted history of $t^{(3)}$. By removing the most basal caterpillar branch b_3 of $t^{(3)}$, we reduce to the 1-rooted history h'' of $t^{(2)}$ (Fig. 5B). Next, by removing the caterpillar branch b_2 of $t^{(2)}$, we reduce to the 2-rooted history h' of $t^{(1)}$ (Fig. 5C). By removing the remaining caterpillar branch b_1 from $t^{(1)}$, we reduce to the 2-rooted history h of $t = t^{(0)}$ (Fig. 5D). Therefore, by the definition of Ω , we have $h' \in \Omega(h)$, $h'' \in \Omega(h')$, and $h''' \in \Omega(h'')$.

By definition, Ω has the property that for each rooted history $h' \in H_{n+1}(t)$, with $n \geq 0$, there exists exactly one rooted history $h \in H_n(t)$ such that $h' \in \Omega(h)$. In other words, for each $n \geq 0$, the set of rooted histories $H_{n+1}(t)$ can be partitioned as a disjoint union,

$$H_{n+1}(t) = \bigsqcup_{h \in H_n(t)} \Omega(h). \quad (11)$$

The set $H_{n+1}(t)$ is therefore generated without double occurrences of any rooted history by applying Ω to the rooted histories in $H_n(t)$. It follows immediately that in performing n iterations of Ω to obtain $\Omega[\dots[\Omega[\Omega(H_0)]]\dots]$ from the set H_0 of rooted histories of $t^{(0)}$, all the rooted histories of $t^{(n)}$ are generated exactly once.

3.2.4 Labels for rooted histories

The operator Ω , starting from the rooted histories of $t^{(n)}$, generates the rooted histories of $t^{(n+1)}$. In this section, we introduce a labeling scheme, giving each m -rooted history h of $t^{(n)}$ a label $L(h) = (n, m)$. We then describe

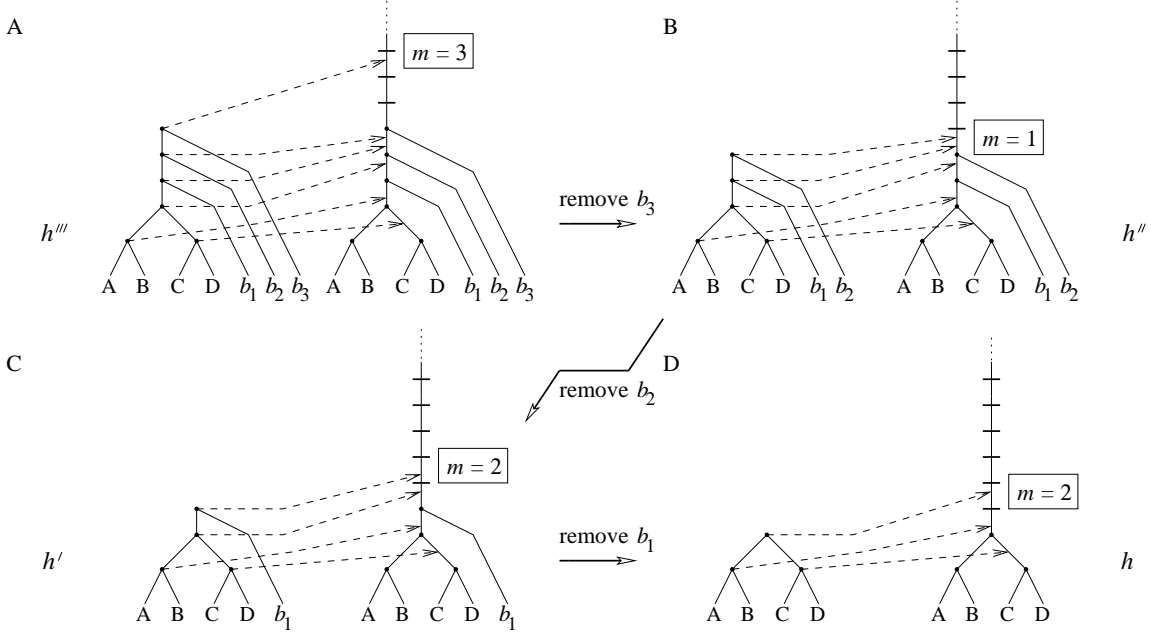


Figure 5: The relationships among rooted histories for sequential members of caterpillar-like families. For a rooted history h''' of $t^{(3)}$, with $t = ((A, B), (C, D))$, the figure sequentially removes caterpillar branches. By definition, a rooted history h' of $t^{(n+1)}$ belongs to the set $\Omega(h)$ if, by removing the most basal caterpillar branch b_{n+1} in $t^{(n+1)}$, we recover the rooted history h of $t^{(n)}$. Note that when we remove the basal caterpillar branch b_{n+1} from $t^{(n+1)}$, the root of $t^{(n+1)}$ —to which the branch b_{n+1} is attached—becomes the boundary between the first and second components of the root-branch of $t^{(n)}$, and is depicted as a horizontal segment. (A) $h''' \in \Omega(h'')$. (B) $h'' \in \Omega(h')$. (C) $h' \in \Omega(h)$. (D) h . For each rooted history, the value of the parameter m , representing the component of the root-branch that receives the image of the root, is shown.

how Ω acts on the labels of the rooted histories, characterizing the set of labels $L[\Omega(h)] = \{L(h') : h' \in \Omega(h)\}$. Our goal is to represent each set H_n of rooted histories of $t^{(n)}$ by the multiset of its labels, reducing the enumeration of $|H_{n,m}|$ to the problem of counting certain ordered pairs (n, m) iteratively generated by simple rules that reflect how the rooted histories in H_{n+1} are generated according to rule Ω from the rooted histories in H_n by eq. (11).

In our labeling scheme, each rooted history $h \in H_n(t)$ that maps the root of $t^{(n)}$ onto the m th component of the root-branch of $t^{(n)}$ receives label $L(h) = (n, m)$. The enumeration of $h_n = |H_{n,1}|$ then reduces to the enumeration of those rooted histories labeled by $(n, 1)$.

Note that a label (n, m) does not uniquely specify an m -rooted history of $t^{(n)}$: a tree $t^{(n)}$ has in general many m -rooted histories, each receiving the label (n, m) . In other words, if $h, \bar{h} \in H_n(t)$ and $L(h) = L(\bar{h})$, then h and \bar{h} are not necessarily the same rooted history of $t^{(n)}$. We will, however, consider for $n \geq 0$ *multisets* of labels in which we find a copy of the label (n, m) for each m -rooted history of $t^{(n)}$.

To characterize how the operator Ω acts on the labels for rooted histories, consider an m -rooted history $h \in H_n(t)$, so that h maps the root of $t^{(n)}$ onto the m th component of the root-branch of $t^{(n)}$. This history is labeled $L(h) = (n, m)$. For instance, taking the seed tree $t = ((A, B), (C, D))$, the history h of $t = t^{(0)}$ depicted in Figure 6A is labeled $L(h) = (0, 3)$, whereas the history h of $t^{(1)}$ in Figure 6C is labeled $L(h) = (1, 1)$.

By applying Ω to a history h of $t^{(n)}$ with $L(h) = (n, m)$, we produce a set of rooted histories $\Omega(h) \subseteq H_{n+1}(t)$. The set of labels for $\Omega(h)$,

$$L[\Omega(h)] = \{L(h') : h' \in \Omega(h)\},$$

is determined according to the rule:

$$L[\Omega(h)] = \begin{cases} \{(n+1, m') : m' \geq m\} & \text{if } m = 1 \\ \{(n+1, m') : m' \geq m-1\} & \text{if } m \geq 2, \end{cases} \quad (12)$$

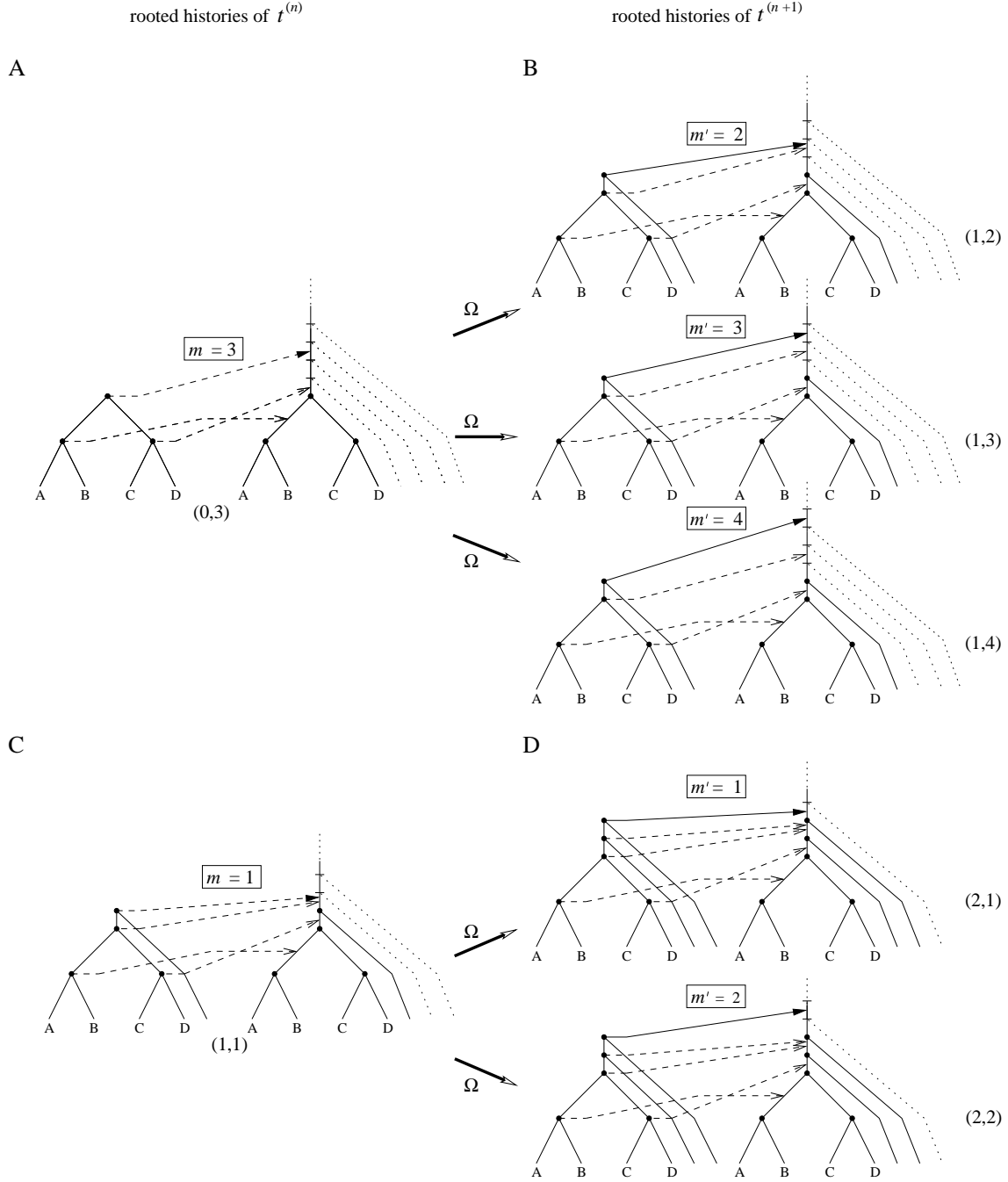


Figure 6: Generation of rooted histories of $t^{(n+1)}$ from rooted histories of $t^{(n)}$, as given by rule Ω applied to seed tree $t = ((A, B), (C, D))$. To obtain rooted histories of $t^{(n+1)}$ (right) from rooted histories of $t^{(n)}$ (left), we choose the component m' of the root-branch of $t^{(n+1)}$ onto which the root of $t^{(n+1)}$ is mapped (solid arrows). The smallest among infinitely many possible choices are depicted. For all nodes of $t^{(n+1)}$ except the root, the rooted history generated for $t^{(n+1)}$ coincides with the generating rooted history of $t^{(n)}$ (dashed arrows). (A) A case with $m \geq 2$. A 2-rooted history h of $t^{(0)}$, labeled $(0, 3)$, is shown. (B) $\Omega(h)$ for h in (A). 2-, 3-, and 4-rooted histories of $t^{(1)}$ belonging to $\Omega(h)$ are shown and are labeled $(1, 2)$, $(1, 3)$, and $(1, 4)$, respectively. Because $m \geq 2$, $m' \geq m - 1$ as in eq. (12). (C) A case with $m = 1$. A 1-rooted history h of $t^{(1)}$, labeled $(1, 1)$, is shown. (D) $\Omega(h)$ for h in (C). 1- and 2-rooted histories of $t^{(2)}$ belonging to $\Omega(h)$ are shown and are labeled $(2, 1)$ and $(2, 2)$, respectively. Because $m = 1$, $m' \geq m$.

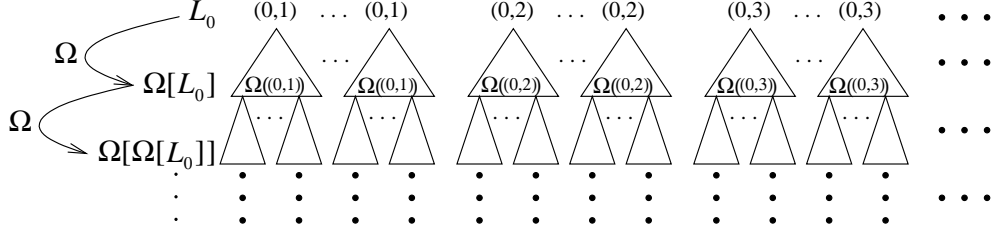


Figure 7: Iterative application of a rule for generating the multiset of the labels of the rooted histories of a tree $t^{(n)}$. The iterative procedure starts with the multiset L_0 that contains those labels of the form $\{(0, m) : m \geq 1\}$ associated with the rooted histories of a seed tree $t = t^{(0)}$. In the first step of the iteration, we apply Ω (eq. (13)) to each label of L_0 . In the second step, we apply Ω to each label resulting from the first step, and so on. The number of m -rooted histories of $t^{(n)}$ corresponds to the number of labels (n, m) , considered with their multiplicity, generated after the n th step of the iteration.

where m' denotes the value of the parameter m —the component of the root-branch of $t^{(n+1)}$ to which the root is mapped—for the rooted histories $h' \in \Omega(h)$ of $t^{(n+1)}$.

The rule in eq. (12) distinguishes between two cases depending on whether the value of the parameter $m = m(h)$ of the generating rooted history h is equal to or exceeds 1. In both cases, the set $L[\Omega(h)]$ contains infinitely many labels, each with its first component equal to $n + 1$, as the labels refer to rooted histories of $t^{(n+1)}$. The value of the second component m' ranges in $[m - 1, \infty)$ if $m \geq 2$, and in $[1, \infty)$ if $m = 1$.

Recall that according to the definition of Ω , from an m -rooted history h of $t^{(n)}$ (Fig. 6A and 6C), we generate an m' -rooted history $h' \in \Omega(h)$ of $t^{(n+1)}$ (Fig. 6B and 6D) by (i) choosing the component m' of the root-branch of $t^{(n+1)}$ onto which h' maps the root of $t^{(n+1)}$, and (ii) letting h' coincide with h on all nodes of $t^{(n+1)}$ except the root. The rooted history h' coincides with h once we remove the most basal caterpillar branch of $t^{(n+1)}$.

Figure 6 illustrates both cases of eq. (12). In step (i), infinitely many choices of m' are possible, because the root-branch of $t^{(n+1)}$ is divided into infinitely many parts. The most basal caterpillar branch in $t^{(n+1)}$ is attached at the border between the first and second components of the root-branch of $t^{(n)}$. Thus, the addition of the $(n + 1)$ st caterpillar branch eliminates a component of the root-branch, so that if the starting rooted history h has $m \geq 2$ (Fig. 6A), then the root of $t^{(n)}$ maps to component $m - 1$ of the root-branch of $t^{(n+1)}$. The root of $t^{(n+1)}$ can map to this same branch, or to any branch m' with $m' \geq m - 1$. For instance, in Figure 6B, one of the rooted histories h' generated by a rooted history h with $m = 3$ has $m' = m - 1 = 2$.

If h has $m = 1$, however, then production of h' is slightly different (Fig. 6C). By definition, the parameter m for a rooted history cannot be smaller than 1. The value $m' = m - 1$ is not permitted in this case, and m' remains greater than or equal to $m = 1$ (Fig. 6D).

3.2.5 From counting rooted histories to counting their labels

The labeling scheme in Section 3.2.4 encodes the application of the operator Ω to the rooted histories of $t^{(n)}$. Now that we have described the set of labels $L[\Omega(h)]$ arising from the label $L(h)$ according to the rule in eq. (12), the problem of counting a set of rooted histories becomes a problem of counting the set of the associated labels along with their multiplicities—or the *multiset* of the labels.

For $n \geq 0$ and $m \geq 1$, we use $\Omega((n, m))$ to denote, with an abuse of notation, the set of labels $L[\Omega(h)]$ when $L(h) = (n, m)$. Recalling that iterative application of Ω to the rooted histories H_0 of tree $t^{(0)}$ generates the rooted histories H_n of $t^{(n)}$, the enumeration of $|H_{n,m}|$ for tree $t = t^{(0)}$ becomes a problem of counting those labels of the form (n, m) that are generated when we iteratively apply the operator Ω as $\Omega[\dots[\Omega[\Omega(L_0)]]\dots]$ starting from the multiset of labels $L_0 = \{L(h) : h \in H_0(t)\}$ (Fig. 7).

Eq. (12) characterizes the set of labels $L[\Omega(h)]$ of the rooted histories in $\Omega(h)$ in terms of the label $L(h)$ of rooted history h . If $L(h) = (n, m)$, then $\Omega((n, m))$ denotes the set of labels $L[\Omega(h)]$. Thus, converting the

notation from histories to labels, eq. (12) becomes

$$\Omega((n, m)) = \begin{cases} \{(n+1, m') : m' \geq m\} & \text{if } m = 1 \\ \{(n+1, m') : m' \geq m-1\} & \text{if } m \geq 2. \end{cases} \quad (13)$$

For the seed tree t , we count $h_{n,m} = |H_{n,m}|$ by evaluating the number of occurrences of the ordered pair (n, m) in the multiset L_n defined as

$$L_n = L[H_n(t)] = \{L(h) : h \in H_n(t)\}. \quad (14)$$

In symbols, we have

$$h_{n,m} = |\{\ell \in L_n : \ell = (n, m)\}|. \quad (15)$$

By eq. (11), each multiset L_n is generated iteratively (Fig. 7). We start with the multiset of labels

$$L_0 = \{L(h) : h \in H_0(t)\}. \quad (16)$$

For each $n \geq 0$, the multiset L_{n+1} is obtained as the union

$$L_{n+1} = \biguplus_{(n,m) \in L_n} \Omega((n, m)), \quad (17)$$

where the symbol \biguplus denotes the union operator for multisets. Thus, in $M = M_1 \biguplus M_2$, if an element x appears n_1 times in M_1 and n_2 times in M_2 , then it appears $n_1 + n_2$ times in M . Eq. (17) provides an iterative generation of the labels for the rooted histories of $H_{n+1}(t)$ from the labels of the rooted histories of $H_n(t)$, retaining information about the multiplicity of occurrences of each label.

3.3 Counting rooted histories with generating functions

We have now obtained eq. (15), which gives an equivalence between the number of m -rooted histories of $t^{(n)}$ and the number of labels (n, m) in the multiset L_n , and eqs. (16) and (17), which give through Ω (eq. (13)) an iterative procedure that generates the family of multisets $(L_n)_{n \geq 0}$. In this section, we translate the iterative procedure into algebraic terms, determining the generating function associated with the integer sequence $(h_n)_{n \geq 0}$.

First, in Section 3.3.1, we characterize a generating function $g(y)$ for the sequence $(h_{0,m})_{m \geq 1}$. Next, in Section 3.3.2, we deduce an equation satisfied by the bivariate generating function $F(y, z)$ for $(h_{n,m})_{n \geq 0, m \geq 1}$. In Section 3.3.3, we solve the equation, obtaining the desired generating function $f(z)$ for the sequence $(h_{n,1})_{n \geq 0}$. This generating function can be written in turn as a function of $g(y)$.

3.3.1 Generating function for the sequence $(h_{0,m})_{m \geq 1}$

In this section, we characterize the generating function $g(y)$ that counts for a given seed tree t the labels in the multiset L_0 describing the labels of the rooted histories of t .

Fix the seed tree t . Recalling the equivalence in eq. (15), define the generating function

$$g(y) = \sum_{(0,m) \in L_0} y^m = \sum_{m=1}^{\infty} h_{0,m} y^m, \quad (18)$$

the m th coefficient of whose power series expansion provides the number $h_{0,m}$ of labels $(0, m)$ appearing in L_0 . By Proposition 1, $h_{0,m}$ can be expressed as a polynomial in the variable m and can thus be decomposed as a finite linear combination of terms of the form m^k , where k is a non-negative integer. That is, for a certain finite set of non-negative integers with largest element K ,

$$h_{0,m} = \sum_{k=0}^K w_k m^k, \quad (19)$$

where the w_k are constants.

We introduce generating functions g_{m^k} , one for each k from 0 to K , in which the m th coefficient is m^k :

$$g_{m^k}(y) = \sum_{m=1}^{\infty} m^k y^m. \quad (20)$$

Because K is finite, the desired generating function $g(y)$ can be written as a finite linear combination of this new collection of generating functions $g_{m^0}(y), g_{m^1}(y), \dots, g_{m^K}(y)$. More precisely, by substituting in eq. (18) the polynomial in eq. (19) and switching the order of summation, we obtain

$$g(y) = \sum_{k=0}^K w_k g_{m^k}(y). \quad (21)$$

We now state a lemma that characterizes the generating functions $g_{m^k}(y)$.

Lemma 1 *For each non-negative integer k from 0 to K , the generating function $g_{m^k}(y)$ in eq. (20) is rational with denominator $(1 - y)^{k+1}$. That is, $g_{m^k}(y)$ has the form*

$$g_{m^k}(y) = \frac{P(y)}{(1 - y)^{k+1}},$$

where $P(y)$ is a polynomial in y .

Proof. We proceed by induction on k . If $k = 0$, then by eq. (20), $g_{m^0}(y) = 1/(1 - y) - 1 = y/(1 - y)$. Assume the inductive hypothesis for $g_{m^k}(y)$. Applying eq. (20) to $g_{m^{k+1}}(y)$, we can recover $g_{m^{k+1}}(y)$ as

$$g_{m^{k+1}}(y) = y \frac{\partial g_{m^k}(y)}{\partial y}, \quad (22)$$

which by the quotient rule for derivatives is a rational function with denominator $(1 - y)^{k+2}$. \square

The proof of the lemma gives a recursive procedure in eq. (22) to compute the functions $g_{m^k}(y)$. By eq. (21), we immediately obtain from the lemma a result about the generating function $g(y)$.

Proposition 2 *The generating function $g(y)$ whose m th coefficient $[y^m]g(y)$ is the number of m -rooted histories $h_{0,m}$ of a seed tree t can be written as a finite linear combination*

$$g(y) = \sum_{j=1}^J q_j \frac{y^{a_j}}{(1 - y)^b}, \quad (23)$$

where $b \geq 1$ and $J \geq 1$ are positive integers, each a_j is a non-negative integer, and the q_j are constants.

As an example, we show how the procedure in Proposition 2 can be applied to determine the generating function $g(y)$ for $t = ((A, B), (C, D))$, the same example seed tree for which we computed the polynomial $h_{0,m}$ via Proposition 1. Recall from Section 3.2.2 that $h_{0,m} = m^2 + 2m + 1$. To obtain the generating function $g(y)$ that has coefficients $[y^m]g(y) = m^2 + 2m + 1$, we sum generating functions for the monomials m^2 , $2m$, and 1. We already know $g_{m^0}(y)$, and by applying eq. (22), we have

$$\begin{aligned} g_{m^0}(y) &= \frac{y}{1 - y} \\ g_{m^1}(y) &= y \frac{\partial g_{m^0}(y)}{\partial y} = \frac{y}{(1 - y)^2} \\ g_{m^2}(y) &= y \frac{\partial g_{m^1}(y)}{\partial y} = \frac{y(y + 1)}{(1 - y)^3}. \end{aligned}$$

Thus,

$$g(y) = g_{m^0}(y) + 2g_{m^1}(y) + g_{m^2}(y) = \frac{y^3 - 3y^2 + 4y}{(1-y)^3}. \quad (24)$$

In eq. (24), $g(y)$ is written as in eq. (23), taking $b = 3$, $J = 3$, $(a_1, a_2, a_3) = (1, 2, 3)$, and $(q_1, q_2, q_3) = (4, -3, 1)$.

3.3.2 Bivariate generating function for the integers $(h_{n,m})_{n \geq 0, m \geq 1}$

Given t , the polynomial nature of $h_{0,m}$ in m enabled us to obtain a generating function for $h_{0,m}$. We now use the iterative procedure in eq. (17) to determine an equation that characterizes the bivariate generating function with coefficients $h_{n,m}$. We represent each label of the form (n, m) by a symbolic algebraic expression in the variables y and z , so that (n, m) is replaced by $z^n y^m$. Let $L = \cup_{n=0}^{\infty} L_n$ be the multiset of all m -rooted histories for all trees $t^{(n)}$. Considering y and z as complex variables in two sufficiently small neighborhoods of 0, we aim to characterize the bivariate function $F(y, z)$ that admits the expansion

$$F(y, z) = \sum_{(n,m) \in L} z^n y^m,$$

where the sum is over all labels in the multiset L and thus has a term for each m -rooted history of each $t^{(n)}$. In particular, the function $F(y, z)$ is the bivariate generating function of the integers $h_{n,m}$, and its Taylor expansion can be written as

$$F(y, z) = \sum_{m=1}^{\infty} \sum_{n=0}^{\infty} h_{n,m} z^n y^m, \quad (25)$$

where the coefficients $h_{n,m}$ appear explicitly.

By differentiating $F(y, z)$ with respect to y and then taking $y = 0$, we obtain

$$\frac{\partial F}{\partial y}(0, z) = \sum_{n=0}^{\infty} h_{n,1} z^n. \quad (26)$$

Thus, for each $n \geq 0$, we have

$$h_n = h_{n,1} = [z^n] \left(\frac{\partial F}{\partial y}(0, z) \right).$$

By representing each label of the form (n, m) by the symbolic expression $z^n y^m$ and assuming the complex variables y and z are sufficiently close to 0, the recursive generation in eq. (17) of the multisets of labels L_0, L_1, L_2, \dots determines an equation for $F(y, z)$, demonstrated in Appendix 1:

$$F(y, z) \left[1 - \frac{z}{y(1-y)} \right] = g(y) - z \frac{\partial F}{\partial y}(0, z). \quad (27)$$

Eq. (27) holds if the complex variables y and z are in two sufficiently small neighborhoods of 0, and it characterizes the generating function $F(y, z)$.

3.3.3 Generating function for the sequence $(h_{n,1})_{n \geq 0}$

We now have an equation satisfied by the bivariate generating function $F(y, z)$. Further, we have eq. (26), which demonstrates that the desired generating function for the sequence $(h_n)_{n \geq 0}$ is obtained from $\frac{\partial F}{\partial y}(0, z)$. By applying the *kernel method* [1, 11], we can determine the power series $\frac{\partial F}{\partial y}(0, z)$ from eq. (27).

The idea of the method consists of coupling the two variables (z, y) as $(z, y(z))$ in such a way that two conditions hold. First, (i) substituting $y = y(z)$ cancels the *kernel* of the equation, that is, the factor $1 - z/[y(1-y)]$

on the left-hand side of eq. (27). Second, (ii) for z near 0, the value of $y(z)$ remains in a sufficiently small neighborhood of $y = 0$, so that eq. (27) still holds near $z = 0$ after substituting $y = y(z)$. This condition is required, as the power series expansion in eq. (25) for $F(y, z)$ has been assumed to be valid in a neighborhood of $(y, z) = (0, 0)$, and the derivation of eq. (27) relies on the fact that y and z are sufficiently close to 0. If the two conditions hold, then

$$z \frac{\partial F}{\partial y}(0, z) = g(y(z)),$$

so that $g(y(z))$ must be a power series for $z = 0$, because so must be $z \frac{\partial F}{\partial y}(0, z)$.

The required substitution couples y and z in such a way that $1 - z/[y(1 - y)] = 0$, so that $y(z) = (1 \pm \sqrt{1 - 4z})/2$. To determine whether to take the negative root $y_1(z)$ or the positive root $y_2(z)$, we note that if z is near 0, then $y_1(z)$ approaches 0, so that $y_1(z)$ lies in a neighborhood of $y = 0$ and $g(y_1(z))$ admits a power series expansion for z near 0. For $y_2(z)$, however, if z is near 0, then $y_2(z)$ approaches 1, and thus, $g(y_2(z))$ is not a power series for z near 0 due to the pole of the function $g(y)$ at $y = 1$ (Proposition 2). The only solution satisfying both (i) and (ii) is consequently

$$Y(z) = y_1(z) = \frac{1 - \sqrt{1 - 4z}}{2}, \quad (28)$$

which, with the generating function $C(z)$ of the Catalan numbers as in eq. (3), satisfies $Y(z) = zC(z)$. Substituting $y = Y(z)$ in eq. (27), we have $\frac{\partial F}{\partial y}(0, z) = g(Y(z))/z$, yielding the following result.

Proposition 3 *Fix tree t . Let $g(y)$ be the generating function associated with the polynomial $h_{0,m}$ (eq. (18)). Let $Y(z)$ be as in eq. (28). Then the generating function $f(z) = \sum_{n=0}^{\infty} h_n z^n$ is given by*

$$f(z) = \frac{\partial F}{\partial y}(0, z) = \frac{g(Y(z))}{z} = \frac{g\left(\frac{1 - \sqrt{1 - 4z}}{2}\right)}{z}. \quad (29)$$

The proposition thus determines the generating function $f(z) = g(Y(z))/z$ for the integer sequence describing the number of matching coalescent histories of the species trees in the caterpillar-like family $(t^{(n)})_{n \geq 0}$. The function g depends on the seed tree t , whereas the function $Y(z)$ is fixed in eq. (28) and does not depend on t .

As an example, recall that for $t = ((A, B), (C, D))$, in eq. (24), we have computed the generating function g for the number $h_{0,m}$ of m -rooted histories of $t = t^{(0)}$. By Proposition 3, the generating function for the number h_n of matching coalescent histories of $t^{(n)}$ is

$$f(z) = \sum_{n=0}^{\infty} h_n z^n = \frac{g\left(\frac{1 - \sqrt{1 - 4z}}{2}\right)}{z} = \frac{4(1 - \sqrt{1 - 4z})(3 - z + \sqrt{1 - 4z})}{z(1 + \sqrt{1 - 4z})^3}.$$

Taking the Taylor expansion of f , we obtain

$$f(z) = 4 + 13z + 42z^2 + 138z^3 + 462z^4 + 1573z^5 + 5434z^6 + 19006z^7 + 67184z^8 + \dots \quad (30)$$

The coefficients h_n accord with the enumeration of matching coalescent histories reported in Corollary 3.9 of [12] and Table 3 of [13] for caterpillar-like families with seed tree $t = ((A, B), (C, D))$, except that those results tabulated numbers of coalescent histories by the number of taxa, whereas here, we use the index of the caterpillar-like family. Thus, in this example, the coefficient of z^n gives the number of matching coalescent histories for a tree with $n + 4$ taxa, as $|t| = 4$. Shifting the index in the formula from [12, 13] to agree with our indexing scheme, we obtain $[(5(n + 4) - 12)/(4(n + 4) - 6)]c_{(n+4)-1} = [(5n + 8)/(4n + 10)]c_{n+3}$ for the number of matching coalescent histories of $t^{(n)}$. This formula gives precisely the coefficients in the Taylor expansion in eq. (30).

3.4 Asymptotic behavior of h_n

From Proposition 3, we have the generating function f that counts the number of matching histories of $t^{(n)}$ for a given fixed seed tree t . Applying techniques of analytic combinatorics as introduced in Section 2.3, we can determine the asymptotic behavior of the coefficients of the generating function

$$\tilde{f}(z) = \sum_{n=1}^{\infty} h_{n-1} z^n = z f(z) = g(Y(z)), \quad (31)$$

with $Y(z)$ as in eq. (28). To simplify notation, we work with \tilde{f} instead of f .

First, in Section 3.4.1, we obtain an asymptotic equivalence between h_n and $\beta_t c_n$, where β_t is a constant depending on the seed tree t , and the c_n are the Catalan numbers (eq. (1)). Next, in Section 3.4.2, we produce a general procedure to determine the constants β_t , employing this procedure to obtain the values of β_t for all seed trees t with $|t| \leq 9$. We demonstrate that our values of β_t accord with constant multiples of the Catalan numbers previously obtained according to a different method [13] for seed trees with $|t| \leq 8$.

3.4.1 A general asymptotic result

Recall that given t , Proposition 2 gives a procedure to determine the rational function g in eq. (31). Writing g as the finite linear combination in eq. (23), the values of b , J , and the $(a_j)_{1 \leq j \leq J}$ and $(q_j)_{1 \leq j \leq J}$ can all be computed.

As noted in Section 2.3, the expansion of \tilde{f} at its dominant singularity characterizes the asymptotic behavior of the coefficients h_{n-1} . In Appendix 2, we obtain the expansion of \tilde{f} at the dominant singularity $z = \frac{1}{4}$,

$$\tilde{f}(z) = \alpha_t + \beta_t \left(-\frac{\sqrt{1-4z}}{2} \right) \pm \mathcal{O}(1-4z) \quad (32)$$

$$\sim \alpha_t + \beta_t \left(-\frac{\sqrt{1-4z}}{2} \right), \quad (33)$$

with

$$\alpha_t = \sum_{j=1}^J 2^{b-a_j} q_j \quad (34)$$

$$\beta_t = \sum_{j=1}^J 2^{b+1-a_j} (a_j + b) q_j. \quad (35)$$

Note that in eq. (32), the seed tree affects only the constants α_t and β_t computed in eqs. (34) and (35) from g , as written in the linear combination in eq. (23). Excluding the constant α_t that does not influence the asymptotic behavior of the coefficients, the main term of the expansion of $\tilde{f}(z)$ (eq. (33)) is the product of the constant β_t and the generating function $-\sqrt{1-4z}/2$, whose n th coefficient is the Catalan number c_{n-1} (eq. (4)).

Theorem VI.4 of [7] indicates that under conditions satisfied by \tilde{f} , the asymptotic coefficients of a generating function as $n \rightarrow \infty$ are obtained from the expansion of the function at the dominant singularity; moreover, the error term in the asymptotic coefficients can be computed from the error term in the singular expansion. Applying the theorem to the expansion in eq. (32), we obtain the asymptotic behavior of the coefficients $[z^n] \tilde{f}(z) = h_{n-1}$.

Proposition 4 *For any seed tree t , when $n \rightarrow \infty$, the number h_n of matching coalescent histories for $t^{(n)}$ satisfies*

$$h_{n-1} = [z^n] \tilde{f}(z) \sim \beta_t [z^n] \left(-\frac{\sqrt{1-4z}}{2} \right) \pm \mathcal{O}\left(\frac{4^n}{n^2}\right) = \beta_t c_{n-1} \pm \mathcal{O}\left(\frac{4^n}{n^2}\right), \quad (36)$$

where β_t is a constant that depends on t . The constant β_t is computed in eq. (35) once the function g , which is defined in eq. (18), has been written as the linear combination in eq. (23).

We immediately obtain the following corollary, corresponding to our initial claim in eq. (6).

Corollary 1 *For any seed tree t , there exists a constant $\beta_t > 0$ (eq. (35)) such that when $n \rightarrow \infty$,*

$$h_n \sim \beta_t c_n. \quad (37)$$

Proof. The result follows from Proposition 4 by noting that if $\beta_t > 0$, then

$$\lim_{n \rightarrow \infty} \frac{h_{n-1}}{\beta_t c_{n-1}} = 1 \pm \lim_{n \rightarrow \infty} \frac{\mathcal{O}(4^n/n^2)}{\beta_t c_{n-1}} = 1.$$

Note that we are claiming $\beta_t > 0$. From the definition of β_t as the sum in eq. (35), because the q_j are permitted to be negative, it is not immediately clear that $\beta_t > 0$. Proposition 4 eliminates the possibility that β_t is negative, as h_{n-1} is necessarily positive. To show that $\beta_t \neq 0$, we note that by eq. (36), $\beta_t = 0$ would give

$$h_{n-1} = \mathcal{O}\left(\frac{4^n}{n^2}\right), \quad (38)$$

so that $h_{n-1}/(4^n/n^2)$ would remain bounded by a constant as $n \rightarrow \infty$.

We now apply the lower bound $h_n \geq c_{n+1}$ from eq. (7). By eq. (7), we have

$$\frac{h_{n-1}}{4^n/n^2} \geq \frac{c_n}{4^n/n^2} = \frac{\sqrt{n}}{\sqrt{\pi}} \frac{c_n}{4^n/(n^{3/2}\sqrt{\pi})}.$$

As $n \rightarrow \infty$, $\sqrt{n}/\sqrt{\pi}$ diverges to ∞ , while $c_n/[4^n/(n^{3/2}\sqrt{\pi})]$ converges to 1 by eq. (5). Therefore, the sequence $h_{n-1}/(4^n/n^2)$ must diverge and eq. (38) cannot hold. Thus, $\beta_t \neq 0$. \square

As an example of Corollary 1, consider $t = ((A, B), (C, D))$. By decomposing the function g expressed in eq. (24) as in eq. (23), we have already obtained the parameters b , J , $(a_j)_{1 \leq j \leq J}$, and $(q_j)_{1 \leq j \leq J}$ in Section 3.3.1. Therefore, computing β_t as in eq. (35), we obtain

$$\beta_t = 2^{1+3-1}(1+3)(4) + 2^{1+3-2}(2+3)(-3) + 2^{1+3-3}(3+3)(1) = 80.$$

Eq. (37) then produces $h_n \sim 80c_n$. Note that the limit $h_n \sim \frac{5}{4}c_{n+3}$ produced for this tree from $h_n = [(5n+8)/(4n+10)]c_{n+3}$ in Section 3.3.3 agrees with the limiting result $h_n \sim 80c_n$. Recalling eq. (2),

$$\frac{h_n}{c_n} = \frac{5n+8}{4n+10} \frac{c_{n+3}}{c_n} \sim \frac{5}{4} \frac{\binom{2n+6}{n+3}/(n+3)}{\binom{2n}{n}/(n+1)} \sim \frac{5}{4} 4^3 = 80.$$

3.4.2 Determining β_t from the seed tree t

We have shown in Corollary 1 that the number of matching coalescent histories h_n for the caterpillar-like family $t^{(n)}$ is, for a constant β_t , asymptotic to $\beta_t c_n$. We can now assemble our results to describe a procedure that given a seed tree t with $|t| \geq 2$ determines both the generating function with coefficients h_n and the constant β_t .

- (i) Determine by eq. (9) the polynomial $e_{t,k}$ in $k \geq 0$ that counts the number of k -extended histories of t .
- (ii) Compute from eq. (8) the polynomial in m that counts for $m \geq 1$ the number of m -rooted histories of t .
- (iii) Obtain the generating function $g(y) = \sum_{m=1}^{\infty} h_{0,m} y^m$ with coefficients $h_{0,m}$ by using Proposition 2.
- (iv) Determine the generating function $f(z) = \sum_{n=0}^{\infty} h_n z^n$ with coefficients h_n by applying Proposition 3.

- (v) Write $g(y)$ as a linear combination according to eq. (23), determining the values of b , J , and the a_j and q_j .
- (vi) Compute the asymptotic constant β_t from eq. (35).

We have programmed this procedure in Mathematica; starting from a given seed tree t , our program `CatFamily.nb` can automatically compute for the caterpillar-like family $t^{(n)}$ the generating function with coefficients h_n and the asymptotic constant β_t . Using this program, we have determined the value of β_t for each seed tree with 9 taxa, collecting the results in Table 1.

Recall that Rosenberg [13] reported the asymptotic constant multiples of the Catalan numbers, β_t^* , which represent the asymptotic numbers of coalescent histories for seed trees with up to 8 taxa, indexing the results by the number of taxa m rather than by the index n of the caterpillar-like family. Also recall that for seed tree t , tree $t^{(n)}$ has $m = |t| + n$ taxa (Fig. 1). In the notation of [13], writing $A_{t_m,1}$ as the number of matching coalescent histories in the caterpillar-like tree with seed tree t and $m \geq |t|$ taxa, we have $h_n = A_{t_m,1}$.

By eq. (5), we have the asymptotic equivalence $c_n \sim c_{n+k}/4^k$ for each positive integer k . Therefore,

$$A_{t_m,1} = h_n \sim \beta_t c_n \sim \frac{\beta_t}{4^{|t|-1}} c_{n+|t|-1} = \beta_t^* c_{m-1}, \quad (39)$$

where the asymptotic constant β_t of Corollary 1 is normalized to obtain

$$\beta_t^* = \frac{\beta_t}{4^{|t|-1}}. \quad (40)$$

This computation converts the asymptotic constant multiple β_t of c_n into a corresponding multiple β_t^* of c_{m-1} , as reported in [13] for small trees. Comparing Table 1 with Table 3 of [13], we see that for the cases examined by [13], the values of β_t^* we compute from the associated β_t agree with the values that were previously reported.

4 Conclusions

In this paper, we have solved a problem left open by [13] on determining the number of coalescent histories for gene trees and species trees that have a matching labeled topology and that belong to a generic caterpillar-like family. We have proven that for any seed tree t , the integer sequence $(h_n)_{n \geq 0}$, whose n th element represents the number of matching coalescent histories of the caterpillar-like tree $t^{(n)}$, grows asymptotically as a constant multiple of the Catalan numbers, that is, $h_n \sim \beta_t c_n$, where the constant term $\beta_t > 0$ depends on the shape of the seed tree t . Rosenberg [13] had previously obtained this result for seed trees with at most 8 taxa; here, by using a succession rule for recursive enumeration and then applying techniques of analytic combinatorics, we have not only proven the existence of the constant β_t for seed trees of any size, we have also produced a procedure that computes the constant β_t , as well as the expression for the generating function of the integers $(h_n)_{n \geq 0}$.

The numerical results on the constants β_t extend the empirical observation of [13] that the caterpillar-like families that produce the largest numbers of matching coalescent histories are those whose seed tree has a high level of balance. By extending from seed trees with $|t| \leq 8$ taxa to those with $|t| = 9$, we observe that the constants β_t for the caterpillar-like families with the largest and smallest numbers of matching coalescent histories become further separated, so that for n large, many more coalescent histories exist by which a gene tree can match the species tree for some species trees than for others. For the 9-taxon seed tree with the largest β_t^* , $\beta_t^* \approx 8.12$ compared to $\beta_t = 1$ for the seed tree with the smallest β_t^* . Our procedure for evaluating β_t and β_t^* as a function of the seed tree can now enable further systematic analyses of the correlates of the constants β_t and β_t^* , to facilitate additional explorations of determinants of the numbers of matching coalescent histories.

Nevertheless, although the constants β_t and β_t^* do depend on the seed tree, we have shown that all caterpillar-like families are asymptotically equivalent in their numbers of matching coalescent histories up to a constant factor. Thus, in considering large trees, the many caterpillar branches contribute to the asymptotic growth

Table 1: Asymptotic constants β_t with $h_n \sim \beta_t c_n$, for seed trees t with 9 taxa.

Seed tree t	β_t	β_t^*	Seed tree t	β_t	β_t^*
	65,536	1		128,864	4,027/2,048
	81,920	5/4		166,624	5,207/2,048
	94,208	23/16		197,296	12,331/4,096
	104,448	51/32		224,704	3,511/1,024
	138,240	135/64		308,576	9,643/2,048
	118,784	29/16		262,000	16,375/4,096
	113,408	443/256		250,272	7,821/2,048
	148,480	145/64		339,504	21,219/4,096
	177,664	347/128		417,632	13,051/2,048
	141,312	69/32		326,240	10,195/2,048
	193,536	189/64		464,128	1,813/256
	121,472	949/512		182,912	1,429/512
	157,888	2,467/1,024		243,904	3,811/1,024
	187,776	1,467/512		296,064	2,313/512
	214,720	3,355/1,024		344,512	5,383/1,024
	296,192	1,157/256		487,808	3,811/512
	251,136	981/256		410,112	801/128
	162,560	635/256		214,016	209/64
	219,136	107/32		306,112	4,783/1,024
	268,288	131/32		294,784	2,303/512
	177,664	347/128		425,216	1,661/256
	249,344	487/128		366,720	2,865/512
	353,536	1,381/256		532,224	2,079/256

Values of β_t appear for each of the 46 unlabeled species trees with 9 taxa. For each species tree t , we also provide the constant $\beta_t^* = \beta_t/4^8$ (eq. (40)). Trees are listed in increasing order by rank as defined in Section 2 of [13]. In the left column, each seed tree t belongs to a caterpillar-like family $(\tilde{t}^{(n)})_n$, with $|\tilde{t}| < 9$. In these cases, we recover the values of β_t^* as determined in Table 3 of [13].

behavior of the number of matching coalescent histories—which follows a multiple of the Catalan numbers—and the seed tree contributes only to the constant by which the Catalan numbers are multiplied. From the viewpoint of computational complexity in evaluating gene tree probabilities according to formulas that sum over matching coalescent histories [4], all caterpillar-like families have the same growth pattern up to a constant.

The extent to which other tree families follow the Catalan sequence in their numbers of matching coalescent histories remains unknown, though we have recently found a family, the *lodgpole* family, for which the number of matching coalescent histories grows faster than with a constant multiple of the Catalan numbers [5]. The use of our substantially different approach employing analytic combinatorics opens new methods for theoretical analysis of coalescent histories and can potentially assist in understanding when Catalan-like growth, the rapid growth of the lodgpole family, and intermediate or perhaps still faster growth patterns will apply.

Appendix 1. The equation satisfied by $F(y, z)$

In this appendix, we complete the derivation of eq. (27) satisfied by $F(y, z)$. In the generating function $F(y, z)$ (eq. (25)), each monomial $z^n y^m$ corresponds to a label $(n, m) \in L_n$ that in turn represents an m -rooted history of $t^{(n)}$. Recall that the multisets of labels L_0, L_1, L_2, \dots (eq. (14)) can be iteratively generated according to eq. (17) through the operator Ω defined in eq. (13), starting from the multiset L_0 . Also recall that by considering the multiset of labels $L = \cup_{n=0}^{\infty} L_n$, we can write $F(y, z) = \sum_{(n,m) \in L} z^n y^m$. We use the iterative generation of the family of multisets $(L_n)_{n \geq 0}$ to obtain an equation for F .

By eq. (13), for $n \geq 0$ and $m \geq 2$, for each occurrence in L_n of a label (n, m) , a copy of each label in set

$$\Omega((n, m)) = \{(n+1, m+j) : j \geq -1\}$$

belongs to the multiset L_{n+1} . Thus, in algebraic terms, each time that an expression $z^n y^m$ with $n \geq 0$ and $m \geq 2$ is counted in the generating function F —written $z^n y^m \in F$ in what follows—the terms $z^{n+1} \sum_{j=m-1}^{\infty} y^j$ appear in F as well. Summing over all possible $z^n y^m \in F$ with $n \geq 0$ and $m \geq 2$, we obtain

$$\sum_{z^n y^m \in F : n \geq 0, m \geq 2} \left(z^{n+1} \sum_{j=m-1}^{\infty} y^j \right) = \frac{z}{y} \sum_{z^n y^m \in F : n \geq 0, m \geq 2} \left(z^n y^m \sum_{j=0}^{\infty} y^j \right). \quad (41)$$

Similarly, for $n \geq 0$ and $m = 1$, for each occurrence in L_n of a label $(n, 1)$, a copy of each label in set $\Omega((n, 1)) = \{(n+1, j) : j \geq 1\}$ appears in the multiset L_{n+1} . Thus, for each term $z^n y \in F$, with $n \geq 0$, the terms $z^{n+1} \sum_{j=1}^{\infty} y^j$ are counted in F as well. Summing these terms for all $z^n y \in F$ with $n \geq 0$, we obtain

$$\sum_{z^n y \in F : n \geq 0} \left(z^{n+1} \sum_{j=1}^{\infty} y^j \right) = zy \sum_{z^n y \in F : n \geq 0} \left(z^n \sum_{j=0}^{\infty} y^j \right). \quad (42)$$

Notice that the sum of the expressions in eqs. (41) and (42) is the algebraic representation of the multiset of labels $L \setminus L_0$. More precisely, each term $z^n y^m \in F$ associated with a label $(n, m) \in L_n$, with $n \geq 1$, is counted—and counted exactly once—in the sum of eqs. (41) and (42). Therefore, to complete the description of F , we require only those terms $z^0 y^m$ associated with labels $(0, m) \in L_0$. These terms are represented by

$$\sum_{(0,m) \in L_0} z^0 y^m = \sum_{m=1}^{\infty} h_{0,m} y^m = g(y), \quad (43)$$

considering that $h_{0,m} = |\{\ell \in L_0 : \ell = (0, m)\}|$ (eq. (15)) and that by definition, $g(y) = \sum_{m=1}^{\infty} h_{0,m} y^m$ (eq. (18)).

We can now equate the full generating function $F(y, z)$ to the sum of eqs. (43), (41), and (42), obtaining

$$F(y, z) = g(y) + \frac{z}{y} \sum_{z^n y^m \in F: n \geq 0, m \geq 2} \left(z^n y^m \sum_{j=0}^{\infty} y^j \right) + zy \sum_{z^n y \in F: n \geq 0} \left(z^n \sum_{j=0}^{\infty} y^j \right). \quad (44)$$

Applying the fact that $\sum_{j=0}^{\infty} y^j = 1/(1-y)$ for y near 0 in the complex plane, we then have

$$F(y, z) = g(y) + \frac{z}{y(1-y)} \left(\sum_{z^n y^m \in F: n \geq 0, m \geq 2} z^n y^m \right) + \frac{zy}{1-y} \left(\sum_{z^n y \in F: n \geq 0} z^n \right). \quad (45)$$

By eq. (25) and the fact that the multisets L_n of labels (n, m) for m -rooted histories of $t^{(n)}$ have $h_{n,m}$ elements,

$$\begin{aligned} \sum_{z^n y \in F: n \geq 0} z^n &= \frac{\partial F}{\partial y}(0, z) \\ \sum_{z^n y^m \in F: n \geq 0, m \geq 2} z^n y^m &= \left(\sum_{z^n y^m \in F: n \geq 0, m \geq 1} z^n y^m \right) - \left(\sum_{z^n y \in F: n \geq 0} z^n y \right) = F(y, z) - y \frac{\partial F}{\partial y}(0, z). \end{aligned}$$

Substituting in eq. (45), the last two expressions yield

$$F(y, z) = g(y) + \frac{z}{y(1-y)} \left(F(y, z) - y \frac{\partial F}{\partial y}(0, z) \right) + \frac{zy}{1-y} \frac{\partial F}{\partial y}(0, z), \quad (46)$$

which can be rewritten as in eq. (27).

Appendix 2. The dominant singularity and singular expansion of $\tilde{f}(z)$

This appendix obtains the singular expansion of $\tilde{f}(z)$ described in eq. (32). In eq. (31), we have defined $\tilde{f}(z)$ as a composition $\tilde{f}(z) = g(Y(z))$, with the internal function $Y(z)$ as in eq. (28) and the external function $g(y)$ as in eq. (23). Owing to the presence of the square root in the expression for $Y(z)$, the dominant singularity of the internal function $Y(z)$ —the singularity nearest the origin of the complex plane—is at $z = \frac{1}{4}$. Computing the value of $Y(z)$ at its dominant singularity, we obtain $Y(\frac{1}{4}) = \frac{1}{2}$. In particular, we have $Y(\frac{1}{4}) < 1$, where 1 is the radius of convergence of the finite series corresponding to the external function g in \tilde{f} . Indeed, it immediately follows from Proposition 2 that $y = 1$ is the dominant singularity of $g(y)$.

As detailed in Section VI.9 of [7], on dominant singularities of compositions, we are in the setting of the subcritical case, in which the inequality $Y(\frac{1}{4}) < 1$ implies that the dominant singularity of $g(Y(z))$ coincides with the dominant singularity $z = \frac{1}{4}$ of the internal function $Y(z)$ rather than the dominant singularity $y = 1$ of the external function $g(y)$. The desired singular expansion of $\tilde{f}(z) = g(Y(z))$ at the dominant singularity $z = \frac{1}{4}$ can be obtained by inserting $y = Y(z)$ in the regular (non-singular) expansion of $g(y)$ at $y = Y(\frac{1}{4}) = \frac{1}{2}$.

To recover the expansion of $g(y)$ at $y = \frac{1}{2}$, we expand and then sum each term $q_j [y^{a_j}/(1-y)^b]$ of the finite linear combination in eq. (23). At $y = \frac{1}{2}$, each of these terms is an analytic function, and we can thus use Taylor's formula to produce the desired expansion. We obtain at $y = \frac{1}{2}$

$$q_j \frac{y^{a_j}}{(1-y)^b} = 2^{b-a_j} q_j + 2^{b+1-a_j} (a_j + b) q_j \left(y - \frac{1}{2} \right) \pm \mathcal{O} \left(\left(y - \frac{1}{2} \right)^2 \right).$$

By summing over the indices $1 \leq j \leq J$ of eq. (23), the expansion of $g(y)$ at $y = \frac{1}{2}$ is

$$g(y) = \alpha_t + \beta_t \left(y - \frac{1}{2} \right) \pm \mathcal{O} \left(\left(y - \frac{1}{2} \right)^2 \right), \quad (47)$$

with the constants α_t and β_t defined as in eqs. (34) and (35). Plugging $y = Y(z)$ from eq. (28) into eq. (47), we finally obtain the singular expansion of $\tilde{f}(z)$ at $z = \frac{1}{4}$ as in eq. (32).

Acknowledgments

We acknowledge grant support from the National Science Foundation (DBI-1146722). A Mathematica notebook `CatFamily.nb` implementing the procedure in Section 3.4.2 for obtaining from a seed tree t the generating function $f(z)$, the coefficients h_n , and the constant β_t is available from the authors.

References

- [1] BANDERIER, C., BOUSQUET-MÉLOU, M., DENISE, A., FLAJOLET, P., GARDY, D., AND GOUYOU-BEAUCHAMPS, D. Generating functions for generating trees. *Discr. Math.* 246 (2002), 29–55.
- [2] BARCUCCI, E., DEL LUNGO, A., PEROGOLA, E., AND PINZANI, R. ECO: a methodology for the enumeration of combinatorial objects. *J. Differ. Equ. Appl.* 5 (1999), 435–490.
- [3] DEGNAN, J. H. *Gene tree distributions under the coalescent process*. PhD thesis, University of New Mexico, Albuquerque, 2005.
- [4] DEGNAN, J. H., AND SALTER, L. A. Gene tree distributions under the coalescent process. *Evolution* 59 (2005), 24–37.
- [5] DISANTO, F., AND ROSENBERG, N. A. Coalescent histories for lodgepole species trees. *Submitted* (2015), xx–xx.
- [6] DUTHEIL, J. Y., GANAPATHY, G., HOBOLTH, A., MAILUND, T., UYENOYAMA, M. K., AND SCHIERUP, M. H. Ancestral population genomics: the coalescent hidden Markov model approach. *Genetics* 183 (2009), 259–274.
- [7] FLAJOLET, P., AND SEDGEWICK, R. *Analytic Combinatorics*. Cambridge University Press, Cambridge, 2009.
- [8] GRAHAM, R. L., KNUTH, D. E., AND PATASHNIK, O. *Concrete Mathematics*, 2nd ed. Addison-Wesley, Boston, 2008.
- [9] HOBOLTH, A., CHRISTENSEN, O. F., MAILUND, T., AND SCHIERUP, M. H. Genomic relationships and speciation times of human, chimpanzee, and gorilla inferred from a coalescent hidden Markov model. *PLoS Genet.* 3 (2007), 294–304.
- [10] HOBOLTH, A., DUTHEIL, J. Y., HAWKS, J., SCHIERUP, M. H., AND MAILUND, T. Incomplete lineage sorting patterns among human, chimpanzee, and orangutan suggest recent orangutan speciation and widespread selection. *Genome Res.* 21 (2011), 349–356.
- [11] PRODINGER, H. The kernel method: a collection of examples. *Sém. Lothar. Combin.* 50 (2004), B50f.
- [12] ROSENBERG, N. A. Counting coalescent histories. *J. Comput. Biol.* 14 (2007), 360–377.
- [13] ROSENBERG, N. A. Coalescent histories for caterpillar-like families. *IEEE/ACM Trans. Comp. Biol. Bioinf.* 10 (2013), 1253–1262.
- [14] ROSENBERG, N. A., AND DEGNAN, J. H. Coalescent histories for discordant gene trees and species trees. *Theor. Pop. Biol.* 77 (2010), 145–151.
- [15] STANLEY, R. P. *Enumerative Combinatorics Volume 2*. Cambridge University Press, New York, 1999.
- [16] THAN, C., RUTHS, D., INNAN, H., AND NAKHLEH, L. Confounding factors in HGT detection: statistical error, coalescent effects, and multiple solutions. *J. Comput. Biol.* 14 (2007), 517–535.